

# ANALYSIS AND MODELLING OF EARLY STAGE DIABETES RISK PREDICTION DATASET

- M.SHOZAB HUSSAIN 23100174
- AHMED TAHIR SHEKHANI 23100197
- JAHANZEB RAZA 23100003
- SYED TALAL 23100176
- ALI ASGHAR 23100198

# DESCRIBING THE DATASET

-> THE DATASET CONTAINS A TOTAL OF **8,857** OBSERVATIONS AND **17** FEATURES. THIS DATASET CONTAINS REPORTS OF DIABETES-RELATED SYMPTOMS OF **521** PERSONS. IT INCLUDES DATA ABOUT PEOPLE INCLUDING SYMPTOMS THAT MAY CAUSE DIABETES.

-> EACH ROW IN THE DATASET REPRESENTS THE RECORD OF A SINGLE PATIENT

-> FEATURES ARE : **AGE, GENDER, POLYURIA, POLYDIPSIA, SUDDEN WEIGHT LOSS, WEAKNESS, POLYPHAGIA, GENITAL THRUSH, VISUAL BLURRING, ITCHING, IRRITABILITY, DELAYED HEALING, PARTIAL PARESIS, MUSCLE STIFFNESS, ALOPECIA, OBESITY, CLASS**

# DATASET CLEANING

- -> MISSING VALUES
- -> BINARY VALUES
- -> COLUMN NAMES
- -> AGE GROUPS



# STATISTICAL INFERENCE AND DATA VISUALIZATION

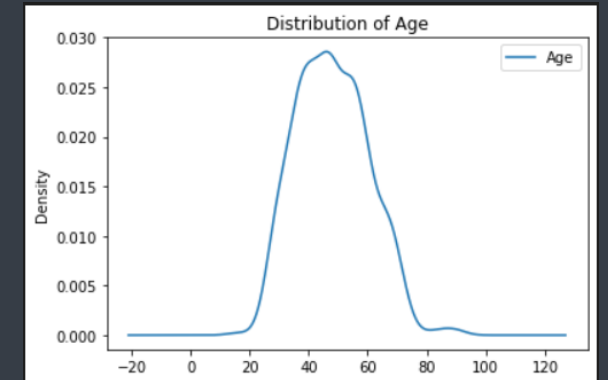
SPREAD OF PATIENTS IN EACH AGE GROUP

NUMBER OF PATIENTS WITH EACH SYMPTOM

CORRELATION MATRIX TO VISUALIZE RELATIONSHIP OF FEATURES

CONCLUSIONS FROM VISUALIZATION:

**DROP ITCHING, DELAYED HEALING AND AGE GROUP**



	Age	Gender	Fatigue	Polydipsia	Insomnia	Weakness	Polyphagia	Gastro-Intestinal	Head-ache	Itching	Intoxication	Delayed healing	Pericarditis	Muscle cramps	Depression	Convulsion	Seizures
Age	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Gender	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Fatigue	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Polydipsia	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Insomnia	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Weakness	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Polyphagia	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Gastro-Intestinal	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Head-ache	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Itching	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Intoxication	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000
Delayed healing	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000	0.000
Pericarditis	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000	0.000
Muscle cramps	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000	0.000
Depression	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	0.000
Convulsion	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000
Seizures	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	1.000



# MACHINE LEARNING MODEL

- PREPARING FEATURE SET AND TEST LABEL
- SPLITTING INTO DATA TEST AND TRAIN DATA
- LOGISTIC REGRESSION
- 5-FOLD CROSS VALIDATION
- PREDICTING TEST LABELS



The Feature Set of our Model:

```
[-0.66, 1. , 0. , ..., 1. , 1. , 1. ],  
[ 0.82, 1. , 0. , ..., 0. , 1. , 0. ],  
[-0.58, 1. , 1. , ..., 1. , 1. , 0. ],  
...,  
[ 0.82, 0. , 1. , ..., 1. , 0. , 1. ],  
[-1.32, 0. , 0. , ..., 0. , 1. , 0. ],  
[-0.5 , 1. , 0. , ..., 0. , 0. , 0. ]
```

Train set: (364, 16) (364, 1)

Test set: (156, 16) (156, 1)

# MODEL RESULTS

Accuracy:

$$\frac{TP+TN}{TP+TN+FP+FN}$$

Recall or Sensitivity:

$$\frac{TP}{TP+FN}$$

Precision:

$$\frac{TP}{TP+FP}$$

F1 Score:

$$\frac{2 * precision * recall}{precision + recall}$$

TP = True positive

TN = True Negative

FP = False Positive

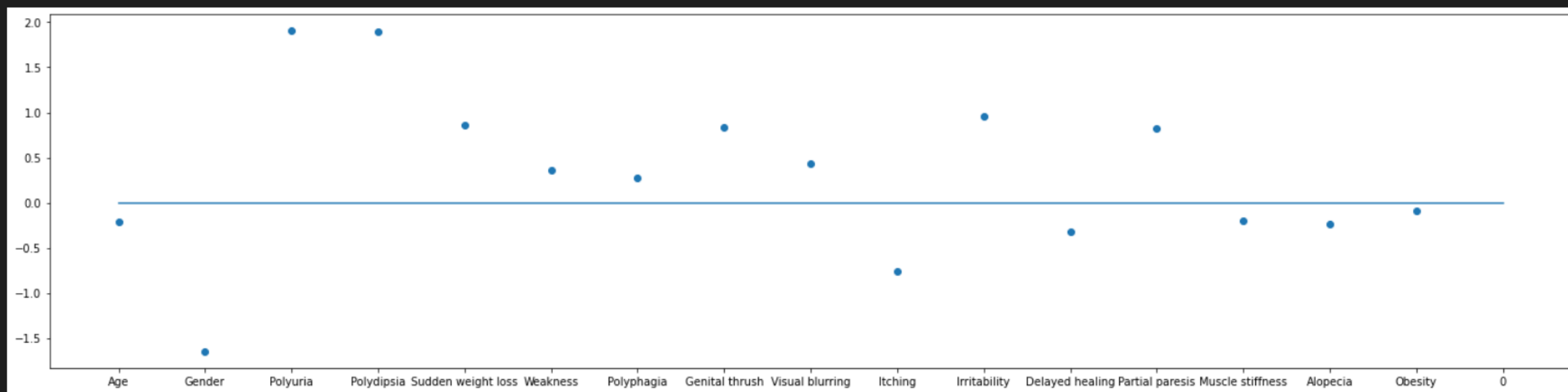
FN = False Negative

Our Dataset result:

Accuracy_train	F1_score_train
0.925824	0.936471
Accuracy_test	F1_score_test
0.942308	0.955665

Precision_train	Recall_train	r2_score_train
0.961353	0.912844	0.691215
Precision_test	Recall_test	r2_score_test
0.960396	0.95098	0.745098

# VISUALIZATION OF WEIGHTS



Weights closer to zero have very less effect on model results (predicted value of either having a diabetes or not)

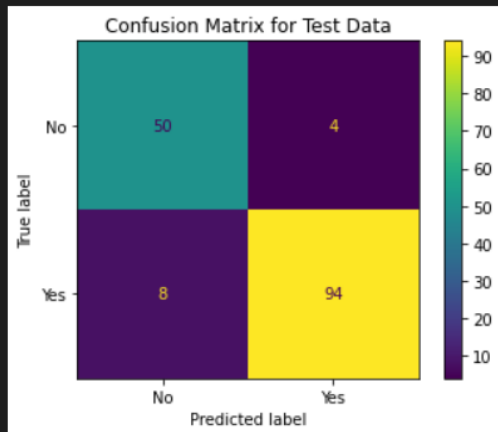
# COMPARING RESULTS BY DROPPING FEATURES

Dropped our proposed features

	Accuracy_train	F1_score_train	Precision_train	Recall_train	r2_score_train
0	0.914835	0.9274	0.947368	0.908257	0.645469

	Accuracy_test	F1_score_test	Precision_test	Recall_test	r2_score_test
0	0.923077	0.94	0.959184	0.921569	0.660131

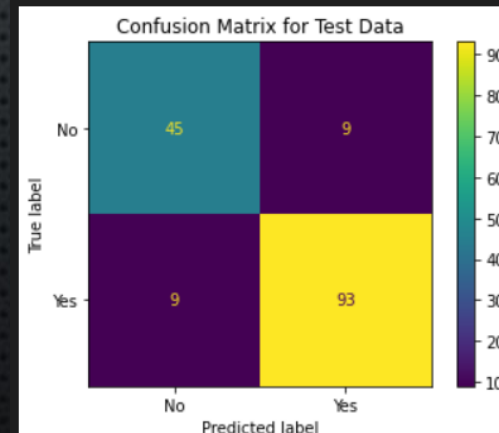


Dropped Gender Column

	Accuracy_train	F1_score_train	Precision_train	Recall_train	r2_score_train
0	0.903846	0.916067	0.959799	0.876147	0.599724

	Accuracy_test	F1_score_test	Precision_test	Recall_test	r2_score_test
0	0.884615	0.911765	0.911765	0.911765	0.490196



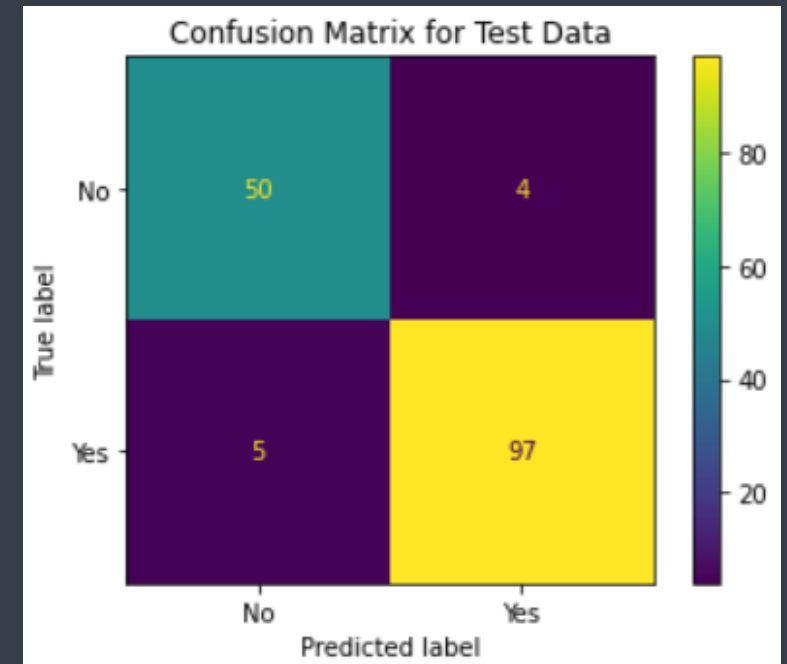




# DID OUR MODEL SOLVE OUR PROBLEM AND ITS REAL-WORLD APPLICATIONS

- THE PROBLEM: CAN OUR MODEL CORRECTLY PREDICT WHETHER A PERSON HAS DIABETES, JUDGING BY THE SIGNS AND SYMPTOMS THEY ARE EXPERIENCING.
- THE OUTCOME: SIGNIFICANTLY SMALL FALSE NEGATIVE PERCENTAGE I.E., APPROX. 3% SUGGESTS THAT IT CAN.

$$\frac{5}{156} \times 100 \approx 3\%$$





# OUR MODEL'S REAL-WORLD APPLICATIONS

- MODEL CAN BE IMPLEMENTED FROM WHERE THE DATA WAS COLLECTED, I.E., BANGLADESH.
- SIMILARLY, OUR MODEL WILL DELIVER PROMISING RESULTS IN COUNTRIES WITH THE SAME DEMOGRAPHIC (E.G., PAKISTAN, INDIA, ETC.).
- QUICK RESULTS AND MINIMAL REQUIREMENTS OFFER CHEAP AND QUICK DIAGNOSIS.
- CAN BE UTILIZED IN AREAS THAT HAVE LIMITED OR NO ACCESS TO HEALTHCARE SERVICES (VILLAGES, TOUGH TERRAIN AREAS, ETC.).